

# Plug-in Approach to Active Learning

Stanislav Minsker<sup>\*,†</sup>

e-mail: [sminsker@math.gatech.edu](mailto:sminsker@math.gatech.edu)

**Abstract:** We present a new active learning algorithm based on non-parametric estimators of the regression function. Our investigation provides probabilistic bounds for the rates of convergence of the generalization error achievable by proposed method over a broad class of underlying distributions. We also prove minimax lower bounds which show that the obtained rates are almost tight.

**Keywords and phrases:** Active learning, selective sampling, model selection, classification, confidence bands.

## 1. Introduction

Let  $(S, \mathcal{B})$  be a measurable space and let  $(X, Y) \in S \times \{-1, 1\}$  be a random couple with unknown distribution  $P$ . The marginal distribution of the design variable  $X$  will be denoted by  $\Pi$ . Let  $\eta(x) := \mathbb{E}(Y|X = x)$  be the regression function. The goal of *binary classification* is to predict label  $Y$  based on the observation  $X$ . Prediction is based on a *classifier* - a measurable function  $f : S \mapsto \{-1, 1\}$ . The quality of a classifier is measured in terms of its generalization error,  $R(f) = \Pr(Y \neq f(X))$ . In practice, the distribution  $P$  remains unknown but the learning algorithm has access to the *training data* - the i.i.d. sample  $(X_i, Y_i)$ ,  $i = 1 \dots n$  from  $P$ . It often happens that the cost of obtaining the training data is associated with labeling the observations  $X_i$  while the pool of observations itself is almost unlimited. This suggests to measure the performance of a learning algorithm in terms of its *label complexity*, the number of labels  $Y_i$  required to obtain a classifier with the desired accuracy. *Active learning* theory is mainly devoted to design and analysis of the algorithms that can take advantage of this modified framework. Most of these procedures can be characterized by the following property: at each step  $k$ , observation  $X_k$  is sampled from a distribution  $\hat{\Pi}_k$  that depends on previously obtained  $(X_i, Y_i)$ ,  $i \leq k - 1$  (while passive learners obtain all available training data at the same time).  $\hat{\Pi}_k$  is designed to be supported on a set where classification is difficult and requires more labeled

---

<sup>\*</sup>Partially supported by ARC Fellowship, NSF Grants DMS-0906880 and CCF-0808863

<sup>†</sup>Mailing address: 686 Cherry street, School of Mathematics, Atlanta, GA 30332-0160

data to be collected. The situation when active learners outperform passive algorithms might occur when the so-called *Tsybakov's low noise assumption* is satisfied: there exist constants  $B, \gamma > 0$  such that

$$\forall t > 0, \Pi(x : |\eta(x)| \leq t) \leq Bt^\gamma \quad (1.1)$$

This assumption provides a convenient way to characterize the noise level of the problem and will play a crucial role in our investigation.

The topic of active learning is widely present in the literature; see Balcan et al. [3], Hanneke [7], Castro and Nowak [4] for review. It was discovered that in some cases the generalization error of a resulting classifier can converge to zero exponentially fast with respect to its label complexity (while the best rate for passive learning is usually polynomial with respect to the cardinality of the training data set). However, available algorithms that adapt to the unknown parameters of the problem ( $\gamma$  in Tsybakov's low noise assumption, regularity of the decision boundary) involve empirical risk minimization with binary loss, along with other computationally hard problems, see Balcan et al. [2], Hanneke [7]. On the other hand, the algorithms that can be effectively implemented, as in Castro and Nowak [4], are not adaptive.

The majority of the previous work in the field was done under standard complexity assumptions on the set of possible classifiers (such as polynomial growth of the covering numbers). Castro and Nowak [4] derived their results under the regularity conditions on the decision boundary and the noise assumption which is slightly more restrictive than (1.1). Essentially, they proved that if the decision boundary is a graph of the Hölder smooth function  $g \in \Sigma(\beta, K, [0, 1]^{d-1})$  (see section 2 for definitions) and the noise assumption is satisfied with  $\gamma > 0$ , then the minimax lower bound for the expected excess risk of the active classifier is of order  $C \cdot N^{-\frac{\beta(1+\gamma)}{2\beta+\gamma(d-1)}}$  and the upper bound is  $C(N/\log N)^{-\frac{\beta(1+\gamma)}{2\beta+\gamma(d-1)}}$ , where  $N$  is the label budget. However, the construction of the classifier that achieves an upper bound assumes  $\beta$  and  $\gamma$  to be known.

In this paper, we consider the problem of active learning under classical nonparametric assumptions on the regression function - namely, we assume that it belongs to a certain Hölder class  $\Sigma(\beta, K, [0, 1]^d)$  and satisfies to the low noise condition (1.1) with some positive  $\gamma$ . In this case, the work of Audibert and Tsybakov [1] showed that plug-in classifiers can attain optimal rates in the *passive* learning framework, namely, that the expected excess risk of a classifier  $\hat{g} = \text{sign } \hat{\eta}$  is bounded above by  $CN^{-\frac{\beta(1+\gamma)}{2\beta+d}}$  (which is the optimal rate), where  $\hat{\eta}$  is the local polynomial estimator of the regression function and  $N$  is the size of the training data set. We were able to partially

extend this claim to the case of active learning: first, we obtain minimax lower bounds for the excess risk of an active classifier in terms of its label complexity. Second, we propose a new algorithm that is based on plug-in classifiers, attains almost optimal rates over a broad class of distributions and possesses adaptivity with respect to  $\beta, \gamma$  (within the certain range of these parameters).

The paper is organized as follows: the next section introduces remaining notations and specifies the main assumptions made throughout the paper. This is followed by a qualitative description of our learning algorithm. The second part of the work contains the statements and proofs of our main results - minimax upper and lower bounds for the excess risk.

## 2. Preliminaries

Our *active learning* framework is governed by the following rules:

1. Observations are sampled sequentially:  $X_k$  is sampled from the modified distribution  $\hat{\Pi}_k$  that depends on  $(X_1, Y_1), \dots, (X_{k-1}, Y_{k-1})$ .
2.  $Y_k$  is sampled from the conditional distribution  $P_{Y|X}(\cdot|X = x)$ . Labels are conditionally independent given the feature vectors  $X_i$ ,  $i \leq n$ .

Usually, the distribution  $\hat{\Pi}_k$  is supported on a set where classification is difficult.

Given the probability measure  $\mathbb{Q}$  on  $S \times \{-1, 1\}$ , we denote the integral with respect to this measure by  $\mathbb{Q}g := \int g d\mathbb{Q}$ . Let  $\mathcal{F}$  be a class of bounded, measurable functions. The risk and the excess risk of  $f \in \mathcal{F}$  with respect to the measure  $\mathbb{Q}$  are defined by

$$\begin{aligned} R_{\mathbb{Q}}(f) &:= \mathbb{Q}\mathcal{I}_{y \neq \text{sign } f(x)} \\ \mathcal{E}_{\mathbb{Q}}(f) &:= R_{\mathbb{Q}}(f) - \inf_{g \in \mathcal{F}} R_{\mathbb{Q}}(g), \end{aligned}$$

where  $\mathcal{I}_{\mathcal{A}}$  is the indicator of event  $\mathcal{A}$ . We will omit the subindex  $\mathbb{Q}$  when the underlying measure is clear from the context. Recall that we denoted the distribution of  $(X, Y)$  by  $P$ . The minimal possible risk with respect to  $P$  is

$$R^* = \inf_{g: S \rightarrow [-1, 1]} \Pr(Y \neq \text{sign } g(X)),$$

where the infimum is taken over all measurable functions. It is well known that it is attained for any  $g$  such that  $\text{sign } g(x) = \text{sign } \eta(x)$   $\Pi$  - a.s. Given  $g \in \mathcal{F}$ ,  $A \in \mathcal{B}$ ,  $\delta > 0$ , define

$$\mathcal{F}_{\infty, A}(g; \delta) := \{f \in \mathcal{F} : \|f - g\|_{\infty, A} \leq \delta\},$$

where  $\|f - g\|_{\infty, A} = \sup_{x \in A} |f(x) - g(x)|$ . For  $A \in \mathcal{B}$ , define the function class

$$\mathcal{F}|_A := \{f|_A, f \in \mathcal{F}\}$$

where  $f|_A(x) := f(x)I_A(x)$ . From now on, we restrict our attention to the case  $S = [0, 1]^d$ . Let  $K > 0$ .

**Definition 2.1.** We say that  $g : \mathbb{R}^d \mapsto \mathbb{R}$  belongs to  $\Sigma(\beta, K, [0, 1]^d)$ , the  $(\beta, K, [0, 1]^d)$  - Hölder class of functions, if  $g$  is  $\lfloor \beta \rfloor$  times continuously differentiable and for all  $x, x_1 \in [0, 1]^d$  satisfies

$$|g(x_1) - T_x(x_1)| \leq K \|x - x_1\|_{\infty}^{\beta},$$

where  $T_x$  is the Taylor polynomial of degree  $\lfloor \beta \rfloor$  of  $g$  at the point  $x$ .

**Definition 2.2.**  $\mathcal{P}(\beta, \gamma)$  is the class of probability distributions on  $[0, 1]^d \times \{-1, +1\}$  with the following properties:

1.  $\forall t > 0, \Pi(x : |\eta(x)| \leq t) \leq Bt^{\gamma};$
2.  $\eta(x) \in \Sigma(\beta, K, [0, 1]^d).$

We do not mention the dependence of  $\mathcal{P}(\beta, \gamma)$  on the fixed constants  $B, K$  explicitly, but this should not cause any uncertainty.

Finally, let us define  $\mathcal{P}_U^*(\beta, \gamma)$  and  $\mathcal{P}_U(\beta, \gamma)$ , the subclasses of  $\mathcal{P}(\beta, \gamma)$ , by imposing two additional assumptions. Along with the formal descriptions of these assumptions, we shall try to provide some motivation behind them. The first deals with the marginal  $\Pi$ . For an integer  $M \geq 1$ , let

$$\mathcal{G}_M := \left\{ \left( \frac{k_1}{M}, \dots, \frac{k_d}{M} \right), k_i = 1 \dots M, i = 1 \dots d \right\}$$

be the regular grid on the unit cube  $[0, 1]^d$  with mesh size  $M^{-1}$ . It naturally defines a partition into a set of  $M^d$  open cubes  $R_i, i = 1 \dots M^d$  with edges of length  $M^{-1}$  and vertices in  $\mathcal{G}_M$ . Below, we consider the nested sequence of grids  $\{\mathcal{G}_{2^m}, m \geq 1\}$  and corresponding dyadic partitions of the unit cube.

**Definition 2.3.** We will say that  $\Pi$  is  $(u_1, u_2)$ -regular with respect to  $\{\mathcal{G}_{2^m}\}$  if for any  $m \geq 1$ , any element of the partition  $R_i, i \leq 2^{dm}$  such that  $R_i \cap \text{supp}(\Pi) \neq \emptyset$ , we have

$$u_1 \cdot 2^{-dm} \leq \Pi(R_i) \leq u_2 \cdot 2^{-dm}. \quad (2.1)$$

where  $0 < u_1 \leq u_2 < \infty$ .

**Assumption 1.**  $\Pi$  is  $(u_1, u_2)$  - regular.

In particular,  $(u_1, u_2)$ -regularity holds for the distribution with a density  $p$  on  $[0, 1]^d$  such that  $0 < u_1 \leq p(x) \leq u_2 < \infty$ .

Let us mention that our definition of regularity is of rather technical nature; for most of the paper, the reader might think of  $\Pi$  as being uniform on  $[0, 1]^d$  (however, we need slightly more complicated marginal to construct the minimax lower bounds for the excess risk). It is known that estimation of regression function in sup-norm is sensitive to the geometry of design distribution, mainly because the quality of estimation depends on the *local* amount of data at every point; conditions similar to our *assumption 1* were used in the previous works where this problem appeared, e.g., *strong density assumption* in Audibert and Tsybakov [1] and *assumption D* in Gaïffas [5]. Another useful characteristic of  $(u_1, u_2)$ -regular distribution  $\Pi$  is that this property is stable with respect to restrictions of  $\Pi$  to certain subsets of its support. This fact fits the active learning framework particularly well.

**Definition 2.4.** We say that  $\mathbb{Q}$  belongs to  $\mathcal{P}_U(\beta, \gamma)$  if  $\mathbb{Q} \in \mathcal{P}(\beta, \gamma)$  and *assumption 1* is satisfied for some  $u_1, u_2$ .

The second assumption is crucial in derivation of the upper bounds. The space of piecewise-constant functions which is used to construct the estimators of  $\eta(x)$  is defined via

$$\mathcal{F}_m = \left\{ \sum_{i=1}^{2^{dm}} \lambda_i I_{R_i}(\cdot) : |\lambda_i| \leq 1, i = 1 \dots 2^{dm} \right\},$$

where  $\{R_i\}_{i=1}^{2^{dm}}$  forms the dyadic partition of the unit cube. Note that  $\mathcal{F}_m$  can be viewed as a  $\|\cdot\|_\infty$ -unit ball in the linear span of first  $2^{dm}$  Haar basis functions in  $[0, 1]^d$ . Moreover,  $\{\mathcal{F}_m, m \geq 1\}$  is a nested family, which is a desirable property for the model selection procedures. By  $\bar{\eta}_m(x)$  we denote the  $L_2(\Pi)$ -projection of the regression function onto  $\mathcal{F}_m$ .

We will say that the set  $A \subset [0, 1]^d$  *approximates the decision boundary*  $\{x : \eta(x) = 0\}$  if there exists  $t > 0$  such that

$$\{x : |\eta(x)| \leq t\}_\Pi \subseteq A_\Pi \subseteq \{x : |\eta(x)| \leq 3t\}_\Pi, \quad (2.2)$$

where for any set  $A$  we define  $A_\Pi := A \cap \text{supp}(\Pi)$ . The most important example we have in mind is the following: let  $\hat{\eta}$  be some estimator of  $\eta$  with  $\|\hat{\eta} - \eta\|_{\infty, \text{supp}(\Pi)} \leq t$ , and define the  $2t$ -band around  $\eta$  by

$$\hat{F} = \left\{ f : \hat{\eta}(x) - 2t \leq f(x) \leq \hat{\eta}(x) + 2t \ \forall x \in [0, 1]^d \right\}$$

Take  $A = \{x : \exists f_1, f_2 \in \hat{F} \text{ s.t. } \text{sign } f_1(x) \neq \text{sign } f_2(x)\}$ , then it is easy to see that  $A$  satisfies (2.2). Modified design distributions used by our algorithm are supported on the sets with similar structure.

Let  $\sigma(\mathcal{F}_m)$  be the sigma-algebra generated by  $\mathcal{F}_m$  and  $A \in \sigma(\mathcal{F}_m)$ .

**Assumption 2.** *There exists  $B_2 > 0$  such that for all  $m \geq 1$ ,  $A \in \sigma(\mathcal{F}_m)$  satisfying (2.2) and such that  $A_\Pi \neq \emptyset$  the following holds true:*

$$\int_{[0,1]^d} (\eta - \bar{\eta}_m)^2 \Pi(dx|x \in A_\Pi) \geq B_2 \|\eta - \bar{\eta}_m\|_{\infty, A_\Pi}^2$$

Appearance of *assumption 2* is motivated by the structure of our learning algorithm - namely, it is based on adaptive confidence bands for the regression function. Nonparametric confidence bands is a big topic in statistical literature, and the review of this subject is not our goal. We just mention that it is impossible to construct adaptive confidence bands of optimal size over the whole  $\bigcup_{\beta \leq 1} \Sigma(\beta, K, [0,1]^d)$ . Hoffmann and Nickl [8], Low [11] discuss the subject in details. However, it is possible to construct adaptive  $L_2$  - confidence balls (see an example following Theorem 6.1 in Koltchinskii [10]). For functions satisfying *assumption 2*, this fact allows to obtain confidence bands of desired size. In particular,

- (a) functions that are differentiable, with gradient being bounded away from 0 in the vicinity of decision boundary;
- (b) Lipschitz continuous functions that are convex in the vicinity of decision boundary

satisfy *assumption 2*. For precise statements, see Propositions A.1, A.2 in Appendix A. A different approach to adaptive confidence bands in case of one-dimensional density estimation is presented in Giné and Nickl [6]. Finally, we define  $\mathcal{P}_U^*(\beta, \gamma)$ :

**Definition 2.5.** *We say that  $\mathbb{Q}$  belongs to  $\mathcal{P}_U^*(\beta, \gamma)$  if  $\mathbb{Q} \in \mathcal{P}_U(\beta, \gamma)$  and *assumption 2* is satisfied for some  $B_2 > 0$ .*

### 2.1. Learning algorithm

Now we give a brief description of the algorithm, since several definitions appear naturally in this context. First, let us emphasize that *the marginal distribution  $\Pi$  is assumed to be known to the learner*. This is not a restriction, since we are not limited in the use of unlabeled data and  $\Pi$  can be

estimated to any desired accuracy. Our construction is based on so-called *plug-in* classifiers of the form  $\hat{f}(\cdot) = \text{sign } \hat{\eta}(\cdot)$ , where  $\hat{\eta}$  is a piecewise-constant estimator of the regression function. As we have already mentioned above, it was shown in Audibert and Tsybakov [1] that in the passive learning framework plug-in classifiers attain optimal rate for the excess risk of order  $N^{-\frac{\beta(1+\gamma)}{2\beta+d}}$ , with  $\hat{\eta}$  being the local polynomial estimator.

Our active learning algorithm iteratively improves the classifier by constructing shrinking confidence bands for the regression function. On every step  $k$ , the piecewise-constant estimator  $\hat{\eta}_k$  is obtained via the model selection procedure which allows adaptation to the unknown smoothness (for Hölder exponent  $\leq 1$ ). The estimator is further used to construct a confidence band  $\hat{\mathcal{F}}_k$  for  $\eta(x)$ . The *active set* associated with  $\hat{\mathcal{F}}_k$  is defined as

$$\hat{A}_k = A(\hat{\mathcal{F}}_k) := \left\{ x \in \text{supp}(\Pi) : \exists f_1, f_2 \in \hat{\mathcal{F}}_k, \text{sign } f_1(x) \neq \text{sign } f_2(x) \right\}$$

Clearly, this is the set where the confidence band crosses zero level and where classification is potentially difficult.  $\hat{A}_k$  serves as a support of the modified distribution  $\hat{\Pi}_{k+1}$ : on step  $k+1$ , label  $Y$  is requested only for observations  $X \in \hat{A}_k$ , forcing the labeled data to concentrate in the domain where higher precision is needed. This allows one to obtain a tighter confidence band for the regression function restricted to the active set. Since  $\hat{A}_k$  approaches the decision boundary, its size is controlled by the low noise assumption. The algorithm does not require a priori knowledge of the noise and regularity parameters, being adaptive for  $\gamma > 0, \beta \leq 1$ . Further details are given in

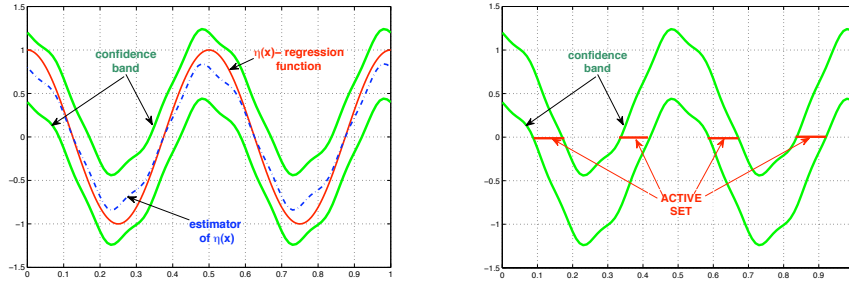


FIG 1. Active Learning Algorithm

section 3.2.

## 2.2. Comparison inequalities

Before proceeding to the main results, let us recall the well-known connections between the binary risk and the  $\|\cdot\|_\infty, \|\cdot\|_{L_2(\Pi)}$  - norm risks:

**Proposition 2.1.** *Under the low noise assumption,*

$$R_P(f) - R^* \leq D_1 \|(f - \eta) \mathcal{I} \{\text{sign } f \neq \text{sign } \eta\}\|_\infty^{1+\gamma}; \quad (2.3)$$

$$R_P(f) - R^* \leq D_2 \|(f - \eta) \mathcal{I} \{\text{sign } f \neq \text{sign } \eta\}\|_{L_2(\Pi)}^{\frac{2(1+\gamma)}{2+\gamma}}; \quad (2.4)$$

$$R_P(f) - R^* \geq D_3 \Pi(\text{sign } f \neq \text{sign } \eta)^{\frac{1+\gamma}{\gamma}} \quad (2.5)$$

*Proof.* For (2.3) and (2.4), see Audibert and Tsybakov [1], lemmas 5.1, 5.2 respectively, and for (2.5)—Koltchinskii [10], lemma 5.2.  $\square$

## 3. Main results

The question we address below is: what are the best possible rates that can be achieved by active algorithms in our framework and how these rates can be attained.

### 3.1. Minimax lower bounds for the excess risk

The goal of this section is to prove that for  $P \in \mathcal{P}(\beta, \gamma)$  no active learner can output a classifier with expected excess risk converging to zero faster than  $N^{-\frac{\beta(1+\gamma)}{2\beta+d-\beta\gamma}}$ . Our result builds upon the minimax bounds of Audibert and Tsybakov [1], Castro and Nowak [4].

**Remark** The theorem below is proved for a smaller class  $\mathcal{P}_U^*(\beta, \gamma)$ , which implies the result for  $\mathcal{P}(\beta, \gamma)$ .

**Theorem 3.1.** *Let  $\beta, \gamma, d$  be such that  $\beta\gamma \leq d$ . Then there exists  $C > 0$  such that for all  $n$  large enough and for any active classifier  $\hat{f}_n(x)$  we have*

$$\sup_{P \in \mathcal{P}_U^*(\beta, \gamma)} \mathbb{E} R_P(\hat{f}_n) - R^* \geq C N^{-\frac{\beta(1+\gamma)}{2\beta+d-\beta\gamma}}$$

*Proof.* We proceed by constructing the appropriate family of classifiers  $f_\sigma(x) = \text{sign } \eta_\sigma(x)$ , in a way similar to Theorem 3.5 in Audibert and Tsybakov [1], and then apply Theorem 2.5 from Tsybakov [13]. We present it below for reader's convenience.



**Theorem 3.2.** Let  $\Sigma$  be a class of models,  $d : \Sigma \times \Sigma \mapsto \mathbb{R}$  - the pseudometric and  $\{P_f, f \in \Sigma\}$  - a collection of probability measures associated with  $\Sigma$ . Assume there exists a subset  $\{f_0, \dots, f_M\}$  of  $\Sigma$  such that

1.  $d(f_i, f_j) \geq 2s > 0 \forall 0 \leq i < j \leq M$
2.  $P_{f_j} \ll P_{f_0}$  for every  $1 \leq j \leq M$
3.  $\frac{1}{M} \sum_{j=1}^M \text{KL}(P_{f_j}, P_{f_0}) \leq \alpha \log M, \quad 0 < \alpha < \frac{1}{8}$

Then

$$\inf_{\hat{f}} \sup_{f \in \Sigma} P_f \left( d(\hat{f}, f) \geq s \right) \geq \frac{\sqrt{M}}{1 + \sqrt{M}} \left( 1 - 2\alpha - \sqrt{\frac{2\alpha}{\log M}} \right)$$

where the infimum is taken over all possible estimators of  $f$  based on a sample from  $P_f$  and  $\text{KL}(\cdot, \cdot)$  is the Kullback-Leibler divergence.

Going back to the proof, let  $q = 2^l$ ,  $l \geq 1$  and

$$G_q := \left\{ \left( \frac{2k_1 - 1}{2q}, \dots, \frac{2k_d - 1}{2q} \right), \quad k_i = 1 \dots q, \quad i = 1 \dots d \right\}$$

be the grid on  $[0, 1]^d$ . For  $x \in [0, 1]^d$ , let

$$n_q(x) = \operatorname{argmin} \{ \|x - x_k\|_2 : x_k \in G_q \}$$

If  $n_q(x)$  is not unique, we choose the one with smallest  $\|\cdot\|_2$  norm. The unit cube is partitioned with respect to  $G_q$  as follows:  $x_1, x_2$  belong to the same subset if  $n_q(x_1) = n_q(x_2)$ . Let  $' \succ'$  be some order on the elements of  $G_q$  such that  $x \succ y$  implies  $\|x\|_2 \geq \|y\|_2$ . Assume that the elements of the partition are enumerated with respect to the order of their centers induced by  $' \succ'$ :

$[0, 1]^d = \bigcup_{i=1}^{q^d} R_i$ . Fix  $1 \leq m \leq q^d$  and let

$$S := \bigcup_{i=1}^m R_i$$

Note that the partition is ordered in such a way that there always exists  $1 \leq k \leq q\sqrt{d}$  with

$$B_+ \left( 0, \frac{k}{q} \right) \subseteq S \subseteq B_+ \left( 0, \frac{k + 3\sqrt{d}}{q} \right), \quad (3.1)$$

where  $B_+(0, R) := \{x \in \mathbb{R}_+^d : \|x\|_2 \leq R\}$ . In other words, (3.1) means that that the difference between the radii of inscribed and circumscribed spherical

sectors of  $S$  is of order  $C(d)q^{-1}$ .

Let  $v > r_1 > r_2$  be three integers satisfying

$$2^{-v} < 2^{-r_1} < 2^{-r_1}\sqrt{d} < 2^{-r_2}\sqrt{d} < 2^{-1} \quad (3.2)$$

Define  $u(x) : \mathbb{R} \mapsto \mathbb{R}_+$  by

$$u(x) := \frac{\int_x^\infty U(t)dt}{\int_{2^{-v}}^{1/2} U(t)dt} \quad (3.3)$$

where

$$U(t) := \begin{cases} \exp\left(-\frac{1}{(1/2-x)(x-2^{-v})}\right), & x \in (2^{-v}, \frac{1}{2}) \\ 0 & \text{else.} \end{cases}$$

Note that  $u(x)$  is an infinitely differentiable function such that  $u(x) = 1$ ,  $x \in [0, 2^{-v}]$  and  $u(x) = 0$ ,  $x \geq \frac{1}{2}$ . Finally, for  $x \in \mathbb{R}^d$  let

$$\Phi(x) := Cu(\|x\|_2)$$

where  $C := C_{L,\beta}$  is chosen such that  $\Phi \in \Sigma(\beta, L, \mathbb{R}^d)$ .

Let  $r_S := \inf \{r > 0 : B_+(0, r) \supseteq S\}$  and

$$A_0 := \left\{ \bigcup_i R_i : R_i \cap B_+(0, r_S + q^{-\frac{\beta\gamma}{d}}) = \emptyset \right\}$$

Note that

$$r_S \leq c \frac{m^{1/d}}{q}, \quad (3.4)$$

since  $\text{Vol}(S) = mq^{-d}$ .

Define  $\mathcal{H}_m = \{P_\sigma : \sigma \in \{-1, 1\}^m\}$  to be the hypercube of probability distributions on  $[0, 1]^d \times \{-1, +1\}$ . The marginal distribution  $\Pi$  of  $X$  is independent of  $\sigma$ : define its density  $p$  by

$$p(x) = \begin{cases} \frac{2^{d(r_1-1)}}{2^{d(r_1-r_2)}-1}, & x \in B_\infty\left(z, \frac{2^{-r_2}}{q}\right) \setminus B_\infty\left(z, \frac{2^{-r_1}}{q}\right), \quad z \in G_q \cap S, \\ c_0, & x \in A_0, \\ 0 & \text{else.} \end{cases}$$

where  $B_\infty(z, r) := \{x : \|x - z\|_\infty \leq r\}$ ,  $c_0 := \frac{1-mq^{-d}}{\text{Vol}(A_0)}$  (note that  $\Pi(R_i) = q^{-d} \quad \forall i \leq m$ ) and  $r_1, r_2$  are defined in (3.2). In particular,  $\Pi$  satisfies

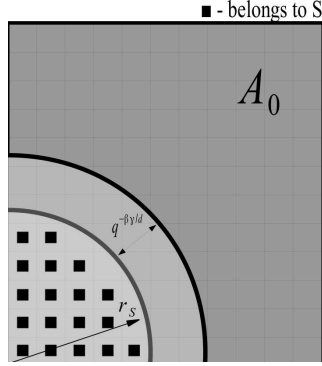


FIG 2. Geometry of the support

assumption 1 since it is supported on the union of dyadic cubes and has bounded above and below on  $\text{supp}(\Pi)$  density.

Let

$$\Psi(x) := u\left(1/2 - q^{\frac{\beta\gamma}{d}} \text{dist}_2(x, B_+(0, r_S))\right),$$

where  $u(\cdot)$  is defined in (3.3) and  $\text{dist}_2(x, A) := \inf \{\|x - y\|_2, y \in A\}$ .

Finally, the regression function  $\eta_\sigma(x) = \mathbb{E}_{P_\sigma}(Y|X = x)$  is defined via

$$\eta_\sigma(x) := \begin{cases} \sigma_i q^{-\beta} \Phi(q[x - n_q(x)]), & x \in R_i, 1 \leq i \leq m \\ \frac{1}{C_{L,\beta}\sqrt{d}} \text{dist}_2(x, B_+(0, r_S))^{\frac{d}{\gamma}} \cdot \Psi(x), & x \in [0, 1]^d \setminus S. \end{cases}$$

The graph of  $\eta_\sigma$  is a surface consisting of small "bumps" spread around  $S$  and tending away from 0 monotonically with respect to  $\text{dist}_2(\cdot, B_+(0, r_S))$  on  $[0, 1]^d \setminus S$ . Clearly,  $\eta_\sigma(x)$  satisfies smoothness requirement, since for  $x \in [0, 1]^d$

$$\text{dist}_2(x, B_+(0, r_S)) = \|x\|_2 - r_S$$

and  $\frac{d}{\gamma} \geq \beta$  by assumption. <sup>1</sup> Let's check that it also satisfies the low noise condition. Since  $|\eta_\sigma| \geq Cq^{-\beta}$  on support of  $\Pi$ , it is enough to consider

---

<sup>1</sup> $\Psi(x)$  can be replaced by 1 unless  $\beta\gamma = d$  and  $\beta$  is an integer, in which case extra smoothness at the boundary of  $B_+(0, r_S)$ , provided by  $\Psi$ , is necessary.

$t = Czq^{-\beta}$  for  $z > 1$ :

$$\begin{aligned} \Pi(|\eta_\sigma(x)| \leq Czq^{-\beta}) &\leq mq^{-d} + \Pi\left(\text{dist}_2(x, B_+(0, r_S)) \leq Cz^{\gamma/d}q^{-\frac{\beta\gamma}{d}}\right) \leq \\ &\leq mq^{-d} + C_2 \left(r_S + Cz^{\gamma/d}q^{-\frac{\beta\gamma}{d}}\right)^d \leq \\ &\leq mq^{-d} + C_3mq^{-d} + C_4z^\gamma q^{-\beta\gamma} \leq \\ &\leq \widehat{C}t^\gamma, \end{aligned}$$

if  $mq^{-d} = O(q^{-\beta\gamma})$ . Here, the first inequality follows from considering  $\eta_\sigma$  on  $S$  and  $A_0$  separately, and second inequality follows from (3.4) and direct computation of the sphere volume.

Finally,  $\eta_\sigma$  satisfies *assumption 2* with some  $B_2 := B_2(q)$  since on  $\text{supp}(\Pi)$

$$0 < c_1(q) \leq \|\nabla \eta_\sigma(x)\|_2 \leq c_2(q) < \infty$$

The next step in the proof is to choose the subset of  $\mathcal{H}$  which is “well-separated”: this can be done due to the following fact (see Tsybakov [13], Lemma 2.9):

**Proposition 3.1** (Gilbert-Varshamov). *For  $m \geq 8$ , there exists*

$$\{\sigma_0, \dots, \sigma_M\} \subset \{-1, 1\}^m$$

*such that  $\sigma_0 = \{1, 1, \dots, 1\}$ ,  $\rho(\sigma_i, \sigma_j) \geq \frac{m}{8} \forall 0 \leq i < k \leq M$  and  $M \geq 2^{m/8}$  where  $\rho$  stands for the Hamming distance.*

Let  $\mathcal{H}' := \{P_{\sigma_0}, \dots, P_{\sigma_M}\}$  be chosen such that  $\{\sigma_0, \dots, \sigma_M\}$  satisfies the proposition above. Next, following the proof of Theorems 1 and 3 in Castro and Nowak [4], we note that  $\forall \sigma \in \mathcal{H}'$ ,  $\sigma \neq \sigma_0$

$$\text{KL}(P_{\sigma, N} \| P_{\sigma_0, N}) \leq 8N \max_{x \in [0, 1]} (\eta_\sigma(x) - \eta_{\sigma_0}(x))^2 \leq 32C_{L, \beta}^2 N q^{-2\beta}, \quad (3.5)$$

where  $P_{\sigma, N}$  is the joint distribution of  $(X_i, Y_i)_{i=1}^N$  under hypothesis that the distribution of couple  $(X, Y)$  is  $P_\sigma$ . Let us briefly sketch the derivation of (3.5); see also the proof of Theorem 1 in Castro and Nowak [4]. Denote

$$\begin{aligned} \bar{X}_k &:= (X_1, \dots, X_k), \\ \bar{Y}_k &:= (Y_1, \dots, Y_k) \end{aligned}$$

Then  $dP_{\sigma, N}$  admits the following factorization:

$$dP_{\sigma, N}(\bar{X}_N, \bar{Y}_N) = \prod_{i=1}^N P_\sigma(Y_i | X_i) dP(X_i | \bar{X}_{i-1}, \bar{Y}_{i-1}),$$

where  $dP(X_i|\bar{X}_{i-1}, \bar{Y}_{i-1})$  does not depend on  $\sigma$  but only on the active learning algorithm. As a consequence,

$$\begin{aligned} \text{KL}(P_{\sigma,N} \| P_{\sigma_0,N}) &= \mathbb{E}_{P_{\sigma,N}} \log \frac{dP_{\sigma,N}(\bar{X}_N, \bar{Y}_N)}{dP_{\sigma_0,N}(\bar{X}_N, \bar{Y}_N)} = \mathbb{E}_{P_{\sigma,N}} \log \frac{\prod_{i=1}^N P_{\sigma}(Y_i|X_i)}{\prod_{i=1}^N P_{\sigma_0}(Y_i|X_i)} = \\ &= \sum_{i=1}^N \mathbb{E}_{P_{\sigma,N}} \left[ \mathbb{E}_{P_{\sigma}} \left( \log \frac{P_{\sigma}(Y_i|X_i)}{P_{\sigma_0}(Y_i|X_i)} \middle| X_i \right) \right] \leq \\ &\leq N \max_{x \in [0,1]^d} \mathbb{E}_{P_{\sigma}} \left( \log \frac{P_{\sigma}(Y_1|X_1)}{P_{\sigma_0}(Y_1|X_1)} \middle| X_1 = x \right) \leq \\ &\leq 8N \max_{x \in [0,1]^d} (\eta_{\sigma}(x) - \eta_{\sigma_0}(x))^2, \end{aligned}$$

where the last inequality follows from Lemma 1, Castro and Nowak [4]. Also, note that we have  $\max_{x \in [0,1]^d}$  in our bounds rather than the average over  $x$  that would appear in the passive learning framework.

It remains to choose  $q, m$  in appropriate way: set  $q \simeq \lfloor C_1 N^{\frac{1}{2\beta+d-\beta\gamma}} \rfloor$  and  $m = \lfloor C_2 q^{d-\beta\gamma} \rfloor$  where  $C_1, C_2$  are such that  $q^d \geq m \geq 1$  and  $32C_{L,\beta}^2 N q^{-2\beta} < \frac{m}{64}$  which is possible for  $N$  big enough. In particular,  $m q^{-d} = O(q^{-\beta\gamma})$ . Together with the bound (3.5), this gives

$$\frac{1}{M} \sum_{\sigma \in \mathcal{H}'} \text{KL}(P_{\sigma} \| P_{\sigma_0}) \leq 32C_u^2 N q^{-2\beta} < \frac{m}{8^2} = \frac{1}{8} \log |\mathcal{H}'|,$$

so that conditions of Theorem 3.2 are satisfied. Setting

$$f_{\sigma}(x) := \text{sign } \eta_{\sigma}(x),$$

we finally have  $\forall \sigma_1 \neq \sigma_2 \in \mathcal{H}'$

$$d(f_{\sigma_1}, f_{\sigma_2}) := \Pi(\text{sign } \eta_{\sigma_1}(x) \neq \text{sign } \eta_{\sigma_2}(x)) \geq \frac{m}{8q^d} \geq C_4 N^{-\frac{\beta\gamma}{2\beta+d-\beta\gamma}},$$

where the lower bound just follows by construction of our hypotheses. Since under the low noise assumption  $R_P(\hat{f}_n) - R^* \geq c \Pi(\hat{f}_n \neq \text{sign } \eta)^{\frac{1+\gamma}{\gamma}}$  (see (2.5)), we conclude that

$$\begin{aligned} &\inf_{\hat{f}_N} \sup_{P \in \mathcal{P}_U^*(\beta, \gamma)} \Pr \left( R_P(\hat{f}_n) - R^* \geq C_4 N^{-\frac{\beta(1+\gamma)}{2\beta+d-\beta\gamma}} \right) \geq \\ &\geq \inf_{\hat{f}_N} \sup_{P \in \mathcal{P}_U^*(\beta, \gamma)} \Pr \left( \Pi(\hat{f}_n(x) \neq \text{sign } \eta_P(x)) \geq \frac{C_4}{2} N^{-\frac{\beta\gamma}{2\beta+d-\beta\gamma}} \right) \geq \tau > 0. \end{aligned}$$

□

### 3.2. Upper bounds for the excess risk

Below, we present a new active learning algorithm which is computationally tractable, adaptive with respect to  $\beta, \gamma$  (in a certain range of these parameters) and can be applied in the nonparametric setting. We show that the classifier constructed by the algorithm attains the rates of Theorem 3.1, up to polylogarithmic factor, if  $0 < \beta \leq 1$  and  $\beta\gamma \leq d$  (the last condition covers the most interesting case when the regression function hits or crosses the decision boundary in the interior of the support of  $\Pi$ ; for detailed statement about the connection between the behavior of the regression function near the decision boundary with parameters  $\beta, \gamma$ , see Proposition 3.4 in Audibert and Tsybakov [1]). The problem of adaptation to higher order of smoothness ( $\beta > 1$ ) is still awaiting its complete solution; we address these questions below in our final remarks.

For the purpose of this section, the regularity assumption reads as follows: there exists  $0 < \beta \leq 1$  such that  $\forall x_1, x_2 \in [0, 1]^d$

$$|\eta(x_1) - \eta(x_2)| \leq B_1 \|x_1 - x_2\|_\infty^\beta \quad (3.6)$$

Since we want to be able to construct non-asymptotic confidence bands, some estimates on the size of constants in (3.6) and *assumption 2* are needed. Below, we will additionally assume that

$$\begin{aligned} B_1 &\leq \log N \\ B_2 &\geq \log^{-1} N, \end{aligned}$$

where  $N$  is the label budget. This can be replaced by any known bounds on  $B_1, B_2$ .

Let  $A \in \sigma(\mathcal{F}_m)$  with  $A_\Pi := A \cap \text{supp}(\Pi) \neq \emptyset$ . Define

$$\hat{\Pi}_A(dx) := \Pi(dx|x \in A_\Pi)$$

and  $d_m := \dim \mathcal{F}_m|_{A_\Pi}$ . Next, we introduce a simple estimator of the regression function on the set  $A_\Pi$ . Given the resolution level  $m$  and an iid sample  $(X_i, Y_i)$ ,  $i \leq N$  with  $X_i \sim \hat{\Pi}_A$ , let

$$\hat{\eta}_{m,A}(x) := \sum_{i: R_i \cap A_\Pi \neq \emptyset} \frac{\sum_{j=1}^N Y_j \mathcal{I}_{R_i}(X_j)}{N \cdot \hat{\Pi}_A(R_i)} \mathcal{I}_{R_i}(x) \quad (3.7)$$

Since we assumed that the marginal  $\Pi$  is known, the estimator is well-defined. The following proposition provides the information about concentration of  $\hat{\eta}_m$  around its mean:

**Proposition 3.2.** *For all  $t > 0$ ,*

$$\begin{aligned} \Pr \left( \max_{x \in A_\Pi} |\hat{\eta}_{m,A}(x) - \bar{\eta}_m(x)| \geq t \sqrt{\frac{2^{dm} \Pi(A)}{u_1 N}} \right) &\leq \\ &\leq 2d_m \exp \left( \frac{-t^2}{2(1 + \frac{t}{3} \sqrt{2^{dm} \Pi(A)/u_1 N})} \right), \end{aligned}$$

*Proof.* This is a straightforward application of the Bernstein's inequality to the random variables

$$S_N^i := \sum_{j=1}^N Y_j \mathcal{I}_{R_i}(X_j), \quad i \in \{i : R_i \cap A_\Pi \neq \emptyset\},$$

and the union bound: indeed, note that  $\mathbb{E}(Y \mathcal{I}_{R_i}(X_j))^2 = \hat{\Pi}_A(R_i)$ , so that

$$\Pr \left( \left| S_N^i - N \int_{R_i} \eta d\hat{\Pi}_A \right| \geq t N \hat{\Pi}_A(R_i) \right) \leq 2 \exp \left( -\frac{N \hat{\Pi}_A(R_i) t^2}{2 + 2t/3} \right),$$

and the rest follows by simple algebra using that  $\hat{\Pi}_A(R_i) \geq \frac{u_1}{2^{dm} \Pi(A)}$  by the  $(u_1, u_2)$ -regularity of  $\Pi$ .  $\square$

Given a sequence of hypotheses classes  $\mathcal{G}_m$ ,  $m \geq 1$ , define the index set

$$\mathcal{J}(N) := \left\{ m \in \mathbb{N} : 1 \leq \dim \mathcal{G}_m \leq \frac{N}{\log^2 N} \right\} \quad (3.8)$$

- the set of possible “resolution levels” of an estimator based on  $N$  classified observations (an upper bound corresponds to the fact that we want the estimator to be consistent). When talking about model selection procedures below, we will implicitly assume that the model index is chosen from the corresponding set  $\mathcal{J}$ . The role of  $\mathcal{G}_m$  will be played by  $\mathcal{F}_m|_A$  for appropriately chosen set  $A$ . We are now ready to present the active learning algorithm followed by its detailed analysis (see Table 1).

**Remark** Note that on every iteration, **Algorithm 1a** uses the whole sample to select the resolution level  $\hat{m}_k$  and to build the estimator  $\hat{\eta}_k$ . While being suitable for practical implementation, this is not convenient for theoretical analysis. We will prove the upper bounds for a slightly modified version: namely, on every iteration  $k$  labeled data is divided into two subsamples  $S_{k,1}$  and  $S_{k,2}$  of approximately equal size,  $|S_{k,1}| \simeq |S_{k,2}| \simeq \left\lfloor \frac{1}{2} N_k \cdot \Pi(\hat{A}_k) \right\rfloor$ .

Algorithm 1a
<b>input</b> label budget $N$ ; confidence $\alpha$ ; $\hat{m}_0 = 0, \hat{\mathcal{F}}_0 := \mathcal{F}_{\hat{m}_0}, \hat{\eta}_0 \equiv 0$ ; $LB := N$ ; // label budget $N_0 := 2^{\lfloor \log_2 \sqrt{N} \rfloor}$ ; $s^{(k)}(m, N, \alpha) := s(m, N, \alpha) := m(\log N + \log \frac{1}{\alpha})$ ; $k := 0$ ; <b>while</b> $LB \geq 0$ <b>do</b> $k := k + 1$ ; $N_k := 2N_{k-1}$ ; $\hat{A}_k := \left\{ x \in [0, 1]^d : \exists f_1, f_2 \in \hat{\mathcal{F}}_{k-1}, \text{sign}(f_1(x)) \neq \text{sign}(f_2(x)) \right\}$ ; <b>if</b> $\hat{A}_k \cap \text{supp}(\Pi) = \emptyset$ <b>or</b> $LB < \lfloor N_k \cdot \Pi(\hat{A}_k) \rfloor$ <b>then</b> <b>break; output</b> $\hat{g} := \text{sign } \hat{\eta}_{k-1}$ <b>else</b> <b>for</b> $i = 1 \dots \lfloor N_k \cdot \Pi(\hat{A}_k) \rfloor$ <b>sample i.i.d</b> $(X_i^{(k)}, Y_i^{(k)})$ <b>with</b> $X_i^{(k)} \sim \hat{\Pi}_k := \Pi(dx x \in \hat{A}_k)$ ; <b>end for</b> ; $LB := LB - \lfloor N_k \cdot \Pi(\hat{A}_k) \rfloor$ ; $\hat{P}_k := \frac{1}{\lfloor N_k \cdot \Pi(\hat{A}_k) \rfloor} \sum_i \delta_{X_i^{(k)}, Y_i^{(k)}} // \text{"active" empirical measure}$ $\hat{m}_k := \text{argmin}_{m \geq \hat{m}_{k-1}} \left[ \inf_{f \in \mathcal{F}_m} \hat{P}_k(Y - f(X))^2 + K_1 \frac{2^{dm} \Pi(\hat{A}_k) + s(m - \hat{m}_{k-1}, N, \alpha)}{\lfloor N_k \cdot \Pi(\hat{A}_k) \rfloor} \right]$ $\hat{\eta}_k := \hat{\eta}_{\hat{m}_k, \hat{A}_k} // \text{see (3.7)}$ $\delta_k := \tilde{D} \cdot \log^2 \frac{N}{\alpha} \sqrt{\frac{2^{d\hat{m}_k}}{N_k}}$ ; $\hat{\mathcal{F}}_k := \left\{ f \in \mathcal{F}_{\hat{m}_k} : f _{\hat{A}_k} \in \mathcal{F}_{\infty, \hat{A}_k}(\hat{\eta}_k; \delta_k), f _{[0,1]^d \setminus \hat{A}_k} \equiv \hat{\eta}_{k-1} _{[0,1]^d \setminus \hat{A}_k} \right\}$ ; <b>end;</b>

TABLE 1  
Active Learning Algorithm

Then  $S_{1,k}$  is used to select the resolution level  $\hat{m}_k$  and  $S_{k,2}$  - to construct  $\hat{\eta}_k$ . We will call this modified version **Algorithm 1b**.

As a first step towards the analysis of **Algorithm 1b**, let us prove the useful fact about the general model selection scheme. Given an iid sample  $(X_i, Y_i)$ ,  $i \leq N$ , set  $s_m = m(s + \log \log_2 N)$ ,  $m \geq 1$  and

$$\hat{m} := \hat{m}(s) = \text{argmin}_{m \in \mathcal{J}(N)} \left[ \inf_{f \in \mathcal{F}_m} P_N(Y - f(X))^2 + K_1 \frac{2^{dm} + s_m}{N} \right] \quad (3.9)$$

$$\bar{m} := \min \left\{ m \geq 1 : \inf_{f \in \mathcal{F}_m} \mathbb{E}(f(X) - \eta(X))^2 \leq K_2 \frac{2^{dm}}{N} \right\} \quad (3.10)$$

**Theorem 3.3.** *There exist an absolute constant  $K_1$  big enough such that, with probability  $\geq 1 - e^{-s}$ ,*

$$\hat{m} \leq \bar{m}$$



*Proof.* See Appendix B.  $\square$

Straightforward application of this result immediately yields the following:

**Corollary 3.1.** *Suppose  $\eta(x) \in \Sigma(\beta, L, [0, 1]^d)$ . Then, with probability  $\geq 1 - e^{-s}$ ,*

$$2^{\hat{m}} \leq C_1 \cdot N^{\frac{1}{2\beta+d}}$$

*Proof.* By definition of  $\bar{m}$ , we have

$$\begin{aligned} \bar{m} &\leq 1 + \max \left\{ m : \inf_{f \in \mathcal{F}_m} \mathbb{E}(f(X) - \eta(X))^2 > K_2 \frac{2^{dm}}{N} \right\} \leq \\ &\leq 1 + \max \left\{ m : L^2 2^{-2\beta m} > K_2 \frac{2^{dm}}{N} \right\}, \end{aligned}$$

and the claim follows.  $\square$

With this bound in hand, we are ready to formulate and prove the main result of this section:

**Theorem 3.4.** *Suppose that  $P \in \mathcal{P}_U^*(\beta, \gamma)$  with  $B_1 \leq \log N$ ,  $B_2 \geq \log^{-1} N$  and  $\beta\gamma \leq d$ . Then, with probability  $\geq 1 - 3\alpha$ , the classifier  $\hat{g}$  returned by **Algorithm 1b** with label budget  $N$  satisfies*

$$R_P(\hat{g}) - R^* \leq \text{Const} \cdot N^{-\frac{\beta(1+\gamma)}{2\beta+d-\beta\gamma}} \log^p \frac{N}{\alpha},$$

where  $p \leq \frac{2\beta\gamma(1+\gamma)}{2\beta+d-\beta\gamma}$  and  $B_1, B_2$  are the constants from (3.6) and assumption 2.

#### Remarks

1. Note that when  $\beta\gamma > \frac{d}{3}$ ,  $N^{-\frac{\beta(1+\gamma)}{2\beta+d-\beta\gamma}}$  is a *fast rate*, i.e., faster than  $N^{-\frac{1}{2}}$ ; at the same time, the passive learning rate  $N^{-\frac{\beta(1+\gamma)}{2\beta+d}}$  is guaranteed to be fast only when  $\beta\gamma > \frac{d}{2}$ , see Audibert and Tsybakov [1].
2. For  $\hat{\alpha} \simeq N^{-\frac{\beta(1+\gamma)}{2\beta+d-\beta\gamma}}$  **Algorithm 1b** returns a classifier  $\hat{g}_{\hat{\alpha}}$  that satisfies

$$\mathbb{E} R_P(\hat{g}_{\hat{\alpha}}) - R^* \leq \text{Const} \cdot N^{-\frac{\beta(1+\gamma)}{2\beta+d-\beta\gamma}} \log^p N.$$

This is a direct corollary of Theorem 3.4 and the inequality

$$\mathbb{E}|Z| \leq t + \|Z\|_{\infty} \Pr(|Z| \geq t)$$

*Proof.* Our main goal is to construct high probability bounds for the size of the active sets defined by **Algorithm 1b**. In turn, these bounds depend on the size of the confidence bands for  $\eta(x)$ , and the previous result (Theorem 3.3) is used to obtain the required estimates. Suppose  $L$  is the number of steps performed by the algorithm before termination; clearly,  $L \leq N$ .

Let  $N_k^{\text{act}} := \lfloor N_k \cdot \Pi(\hat{A}_k) \rfloor$  be the number of labels requested on  $k$ -th step of the algorithm: this choice guarantees that the "density" of labeled examples doubles on every step.

Claim: the following bound for the size of the active set holds uniformly for all  $2 \leq k \leq L$  with probability at least  $1 - 2\alpha$ :

$$\Pi(\hat{A}_k) \leq C N_k^{-\frac{\beta\gamma}{2\beta+d}} \left( \log \frac{N}{\alpha} \right)^{2\gamma} \quad (3.11)$$

It is not hard to finish the proof assuming (3.11) is true: indeed, it implies that the number of labels requested on step  $k$  satisfies

$$N_k^{\text{act}} = \lfloor N_k \Pi(\hat{A}_k) \rfloor \leq C \cdot N_k^{\frac{2\beta+d-\beta\gamma}{2\beta+d}} \left( \log \frac{N}{\alpha} \right)^{2\gamma}$$

with probability  $\geq 1 - 2\alpha$ . Since  $\sum_k N_k^{\text{act}} \leq N$ , one easily deduces that on the last iteration  $L$  we have

$$N_L \geq c \left( \frac{N}{\log^{2\gamma}(N/\alpha)} \right)^{\frac{2\beta+d}{2\beta+d-\beta\gamma}} \quad (3.12)$$

To obtain the risk bound of the theorem from here, we apply inequality (2.3)<sup>2</sup> from proposition 2.1:

$$R_P(\hat{g}) - R^* \leq D_1 \|(\hat{\eta}_L - \eta) \cdot \mathcal{I} \{ \text{sign } \hat{\eta}_L \neq \text{sign } \eta \} \|_{\infty}^{1+\gamma} \quad (3.13)$$

It remains to estimate  $\|\hat{\eta}_L - \eta\|_{\infty, \hat{A}_L}$ : we will show below while proving (3.11) that

$$\|\hat{\eta}_L - \eta\|_{\infty, \hat{A}_L} \leq C \cdot N_L^{-\frac{\beta}{2\beta+d}} \log^2 \frac{N}{\alpha}$$

Together with (3.12) and (3.13), it implies the final result.

To finish the proof, it remains to establish (3.11). Recall that  $\bar{\eta}_k$  stands for the  $L_2(\Pi)$  - projection of  $\eta$  onto  $\mathcal{F}_{\hat{\Pi}_k}$ . An important role in the argument is played by the bound on the  $L_2(\Pi_k)$  - norm of the "bias" ( $\bar{\eta}_k - \eta$ ):

---

<sup>2</sup>alternatively, inequality (2.4) can be used but results in a slightly inferior logarithmic factor.

together with *assumption 2*, it allows to estimate  $\|\bar{\eta}_k - \eta\|_{\infty, \hat{A}_k}$ . The required bound follows from the following oracle inequality: there exists an event  $\mathcal{B}$  of probability  $\geq 1 - \alpha$  such that on this event for every  $1 \leq k \leq L$

$$\begin{aligned} \|\bar{\eta}_k - \eta\|_{L_2(\hat{\Pi}_k)}^2 &\leq \inf_{m \geq \hat{m}_{k-1}} \left[ \inf_{f \in \mathcal{F}_m} \|f - \eta\|_{L_2(\hat{\Pi}_k)}^2 + \right. \\ &\quad \left. + K_1 \frac{2^{dm} \Pi(\hat{A}_k) + (m - \hat{m}_{k-1}) \log(N/\alpha)}{N_k \Pi(\hat{A}_k)} \right] \end{aligned} \quad (3.14)$$

In general form, this inequality is given by Theorem 6.1 in Koltchinskii [10] and provides the estimate for  $\|\hat{\eta}_k - \eta\|_{L_2(\hat{\Pi}_k)}$ , so it automatically implies the weaker bound for the bias term only. To deduce (3.14), we use the mentioned general inequality  $L$  times (once for every iteration) and the union bound. The quantity  $2^{dm} \Pi(\hat{A}_k)$  in (3.14) plays the role of the dimension, which is justified below. Let  $k \geq 1$  be fixed. For  $m \geq \hat{m}_{k-1}$ , consider hypothesis classes

$$\mathcal{F}_m|_{\hat{A}_k} := \left\{ f\mathcal{I}_{\hat{A}_k}, f \in \mathcal{F}_m \right\}$$

An obvious but important fact is that for  $P \in \mathcal{P}_U(\beta, \gamma)$ , the dimension of  $\mathcal{F}_m|_{\hat{A}_k}$  is bounded by  $u_1^{-1} \cdot 2^m \Pi(\hat{A}_k)$ : indeed,

$$\Pi(\hat{A}_k) = \sum_{j: R_j \cap \hat{A}_k \neq \emptyset} \Pi(R_j) \geq u_1 2^{-dm} \cdot \# \left\{ j : R_j \cap \hat{A}_k \neq \emptyset \right\},$$

hence

$$\dim \mathcal{F}_m|_{\hat{A}_k} = \# \left\{ j : R_j \cap \hat{A}_k \neq \emptyset \right\} \leq u_1^{-1} \cdot 2^m \Pi(\hat{A}_k). \quad (3.15)$$

Theorem 3.3 applies conditionally on  $\left\{ X_i^{(j)} \right\}_{i=1}^{N_j}$ ,  $j \leq k-1$  with sample of size  $N_k^{\text{act}}$  and  $s = \log(N/\alpha)$ : to apply the theorem, note that, by definition of  $\hat{A}_k$ , it is independent of  $X_i^{(k)}$ ,  $i = 1 \dots N_k^{\text{act}}$ . Arguing as in Corollary 3.1 and using (3.15), we conclude that the following inequality holds with probability  $\geq 1 - \frac{\alpha}{N}$  for every fixed  $k$ :

$$2^{\hat{m}_k} \leq C \cdot N_k^{\frac{1}{2\beta+d}}. \quad (3.16)$$

Let  $\mathcal{E}_1$  be an event of probability  $\geq 1 - \alpha$  such that on this event bound (3.16) holds for every step  $k$ ,  $k \leq L$  and let  $\mathcal{E}_2$  be an event of probability  $\geq 1 - \alpha$  on which inequalities (3.14) are satisfied. Suppose that event  $\mathcal{E}_1 \cap \mathcal{E}_2$

occurs and let  $k_0$  be a fixed arbitrary integer  $2 \leq k_0 \leq L + 1$ . It is enough to assume that  $\hat{A}_{k_0-1}$  is nonempty (otherwise, the bound trivially holds), so that it contains at least one cube with sidelength  $2^{-\hat{m}_{k_0-2}}$  and

$$\Pi(\hat{A}_{k_0-1}) \geq u_1 2^{-d\hat{m}_{k_0-1}} \geq c N_{k_0}^{-\frac{d}{2\beta+d}} \quad (3.17)$$

Consider inequality (3.14) with  $k = k_0 - 1$  and  $2^m \simeq N_{k_0-1}^{\frac{1}{2\beta+d}}$ . By (3.17), we have

$$\|\bar{\eta}_{k_0-1} - \eta\|_{L_2(\hat{\Pi}_{k_0-1})}^2 \leq C N_{k_0-1}^{-\frac{2\beta}{2\beta+d}} \log^2 \frac{N}{\alpha} \quad (3.18)$$

For convenience and brevity, denote  $\Omega := \text{supp}(\Pi)$ . Now *assumption 2* comes into play: it implies, together with (3.18) that

$$C N_{k_0-1}^{-\frac{\beta}{2\beta+d}} \log \frac{N}{\alpha} \geq \|\bar{\eta}_{k_0-1} - \eta\|_{L_2(\hat{\Pi}_{k_0-1})} \geq B_2 \|\bar{\eta}_{k_0-1} - \eta\|_{\infty, \Omega \cap \hat{A}_{k_0-1}} \quad (3.19)$$

To bound

$$\|\hat{\eta}_{k_0-1}(x) - \bar{\eta}_{k_0-1}(x)\|_{\infty, \Omega \cap \hat{A}_{k_0-1}}$$

we apply Proposition 3.2. Recall that  $\hat{m}_{k_0-1}$  depends only on the subsample  $S_{k_0-1,1}$  but not on  $S_{k_0-1,2}$ . Let

$$\mathcal{T}_k := \left\{ \left\{ X_i^{(j)}, Y_i^{(j)} \right\}_{i=1}^{N_j^{\text{act}}}, j \leq k-1; S_{k,1} \right\}$$

be the random vector that defines  $\hat{A}_k$  and resolution level  $\hat{m}_k$ . Note that  $\mathbb{E}(\hat{\eta}_{k_0-1}(x) | \mathcal{T}_{k_0-1}) = \bar{\eta}_{\hat{m}_{k_0-1}}(x) \quad \forall x \text{ a.s.}$

Proposition 3.2 thus implies

$$\begin{aligned} \Pr \left( \max_{x \in \Omega \cap \hat{A}_{k_0-1}} |\hat{\eta}_{k_0-1}(x) - \bar{\eta}_{\hat{m}_{k_0-1}}(x)| \geq K t \sqrt{\frac{2^{d\hat{m}_{k_0-1}}}{N_{k_0-1}}} \middle| \mathcal{T}_{k_0-1} \right) &\leq \\ &\leq N \exp \left( \frac{-t^2}{2(1 + \frac{t}{3} C_3)} \right). \end{aligned}$$

Choosing  $t = c \log(N/\alpha)$  and taking expectation, the inequality (now unconditional) becomes

$$\Pr \left( \max_{x \in \Omega \cap \hat{A}_{k_0-1}} |\hat{\eta}_{\hat{m}_{k_0-1}}(x) - \bar{\eta}_{\hat{m}_{k_0-1}}(x)| \leq K \sqrt{\frac{2^{d\hat{m}_{k_0-1}} \log^2(N/\alpha)}{N_{k_0-1}}} \right) \geq 1 - \alpha \quad (3.20)$$

Let  $\mathcal{E}_3$  be the event on which (3.20) holds true. Combined, the estimates (3.16), (3.19) and (3.20) imply that on  $\mathcal{E}_1 \cap \mathcal{E}_2 \cap \mathcal{E}_3$

$$\begin{aligned} \|\eta - \hat{\eta}_{k_0-1}\|_{\infty, \Omega \cap \hat{A}_{k_0-1}} &\leq \|\eta - \bar{\eta}_{k_0-1}\|_{\infty, \Omega \cap \hat{A}_{k_0-1}} + \|\bar{\eta}_{k_0-1} - \hat{\eta}_{k_0-1}\|_{\infty, \Omega \cap \hat{A}_{k_0-1}} \\ &\leq \frac{C}{B_2} N_{k_0-1}^{-\frac{\beta}{2\beta+d}} \log \frac{N}{\alpha} + K \sqrt{\frac{2^{d\hat{m}_{k_0-1}} \log^2(N/\alpha)}{N_{k_0-1}}} \leq \\ &\leq (K + C) \cdot N_{k_0-1}^{-\frac{\beta}{2\beta+d}} \log^2 \frac{N}{\alpha} \end{aligned} \quad (3.21)$$

where we used the assumption  $B_2 \geq \log^{-1} N$ . Now the width of the confidence band is defined via

$$\delta_k := 2(K + C) \cdot N_{k_0-1}^{-\frac{\beta}{2\beta+d}} \log^2 \frac{N}{\alpha} \quad (3.22)$$

(in particular,  $\tilde{D}$  from **Algorithm 1a** is equal to  $2(K + C)$ ). With the bound (3.21) available, it is straightforward to finish the proof of the claim. Indeed, by (3.22) and the definition of the active set, the necessary condition for  $x \in \Omega \cap \hat{A}_{k_0}$  is

$$|\eta(x)| \leq 3(K + C) \cdot N_{k_0-1}^{-\frac{\beta}{2\beta+d}} \log^2 \frac{N}{\alpha},$$

so that

$$\begin{aligned} \Pi(\hat{A}_{k_0}) &= \Pi(\Omega \cap \hat{A}_{k_0}) \leq \Pi\left(|\eta(x)| \leq 3(K + C) \cdot N_{k_0-1}^{-\frac{\beta}{2\beta+d}} \log^2 \frac{N}{\alpha}\right) \leq \\ &\leq \tilde{B} N_{k_0-1}^{-\frac{\beta\gamma}{2\beta+d}} \log^{2\gamma} \frac{N}{\alpha} \end{aligned}$$

by the low noise assumption. This completes the proof of the claim since  $\Pr(\mathcal{E}_1 \cap \mathcal{E}_2 \cap \mathcal{E}_3) \geq 1 - 3\alpha$ .  $\square$

We conclude this section by discussing running time of the active learning algorithm. Assume that the algorithm has access to the sampling subroutine that, given  $A \subset [0, 1]^d$  with  $\Pi(A) > 0$ , generates i.i.d.  $(X_i, Y_i)$  with  $X_i \sim \Pi(dx|x \in A)$ .

**Proposition 3.3.** *The running time of Algorithm 1a(1b) with label budget  $N$  is*

$$\mathcal{O}(dN \log^2 N).$$

**Remark** In view of Theorem 3.4, the running time required to output a classifier  $\hat{g}$  such that  $R_P(\hat{g}) - R^* \leq \varepsilon$  with probability  $\geq 1 - \alpha$  is

$$\mathcal{O} \left( \left( \frac{1}{\varepsilon} \right)^{\frac{2\beta+d-\beta\gamma}{\beta(1+\gamma)}} \text{poly} \left( \log \frac{1}{\varepsilon\alpha} \right) \right).$$

*Proof.* We will use the notations of Theorem 3.4. Let  $N_k^{\text{act}}$  be the number of labels requested by the algorithm on step  $k$ . The resolution level  $\hat{m}_k$  is always chosen such that  $\hat{A}_k$  is partitioned into at most  $N_k^{\text{act}}$  dyadic cubes, see (3.8). This means that the estimator  $\hat{\eta}_k$  takes at most  $N_k^{\text{act}}$  distinct values. The key observation is that for any  $k$ , the active set  $\hat{A}_{k+1}$  is always represented as the union of a finite number (at most  $N_k^{\text{act}}$ ) of dyadic cubes: to determine if a cube  $R_j \subset \hat{A}_{k+1}$ , it is enough to take a point  $x \in R_j$  and compare  $\text{sign}(\hat{\eta}_k(x) - \delta_k)$  with  $\text{sign}(\hat{\eta}_k(x) + \delta_k)$ :  $R_j \in \hat{A}_{k+1}$  only if the signs are different (so that the confidence band crosses zero level). This can be done in  $\mathcal{O}(N_k^{\text{act}})$  steps.

Next, resolution level  $\hat{m}_k$  can be found in  $\mathcal{O}(N_k^{\text{act}} \log^2 N)$  steps: there are at most  $\log_2 N_k^{\text{act}}$  models to consider; for each  $m$ ,  $\inf_{f \in \mathcal{F}_m} \hat{P}_k(Y - f(X))^2$  is found explicitly and is achieved for the piecewise-constant

$$\hat{f}(x) = \frac{\sum_i Y_i^{(k)} \mathcal{I}_{R_j}(X_i^{(k)})}{\sum_i \mathcal{I}_{R_j}(X_i^{(k)})}, \quad x \in R_j.$$

Sorting of the data required for this computation is done in  $\mathcal{O}(dN_k^{\text{act}} \log N)$  steps for each  $m$ , so the whole  $k$ -th iteration running time is  $\mathcal{O}(dN_k^{\text{act}} \log^2 N)$ . Since  $\sum_k N_k^{\text{act}} \leq N$ , the result follows.  $\square$

#### 4. Conclusion and open problems

We have shown that active learning can significantly improve the quality of a classifier over the passive algorithm for a large class of underlying distributions. Presented method achieves fast rates of convergence for the excess risk, moreover, it is adaptive (in the certain range of smoothness and noise parameters) and involves minimization only with respect to quadratic loss (rather than the 0 – 1 loss).

The natural question related to our results is:

- Can we implement adaptive smooth estimators in the learning algorithm to extend our results beyond the case  $\beta \leq 1$ ?

The answer to this second question is so far an open problem. Our conjecture is that the correct rate of convergence for the excess risk is  $N^{-\frac{\beta(1+\gamma)}{2\beta+d-\gamma(\beta\wedge 1)}}$ , up to logarithmic factors, which coincides with presented results for  $\beta \leq 1$ . This rate can be derived from an argument similar to the proof of Theorem 3.4 under the assumption that on every step  $k$  one could construct an estimator  $\hat{\eta}_k$  with

$$\|\eta - \hat{\eta}_k\|_{\infty, \hat{A}_k} \lesssim N_k^{-\frac{\beta}{2\beta+d}}.$$

At the same time, the active set associated to  $\hat{\eta}_k$  should maintain some structure which is suitable for the iterative nature of the algorithm. Transforming these ideas into a rigorous proof is a goal of our future work.

### Acknowledgements

I want to express my deepest gratitude to my Ph.D. advisor, Dr. Vladimir Koltchinskii, for his support and numerous helpful discussions.

I am grateful to the anonymous reviewers for carefully reading the manuscript. Their insightful and wise suggestions helped to improve the quality of presentation and results.

I would like to acknowledge support for this project from the National Science Foundation (NSF Grants DMS-0906880 and CCF-0808863) and by the Algorithms and Randomness Center, Georgia Institute of Technology, through the ARC Fellowship.

### Appendix A: Functions satisfying assumption 2

In the propositions below, we will assume for simplicity that the marginal distribution  $\Pi$  is absolutely continuous with respect to Lebesgue measure with density  $p(x)$  such that

$$0 < p_1 \leq p(x) \leq p_2 < \infty \text{ for all } x \in [0, 1]^d \quad (\text{A.1})$$

Given  $t \in (0, 1]$ , define  $A_t := \{x : |\eta(x)| \leq t\}$ .

**Proposition A.1.** *Suppose  $\eta$  is Lipschitz continuous with Lipschitz constant  $S$ . Assume also that for some  $t_* > 0$  we have*

- (a)  $\Pi(A_{t_*/3}) > 0$ ;
- (b)  $\eta$  is twice differentiable for all  $x \in A_{t_*}$ ;
- (c)  $\inf_{x \in A_{t_*}} \|\nabla \eta(x)\|_1 \geq s > 0$ ;
- (d)  $\sup_{x \in A_{t_*}} \|D^2 \eta(x)\| \leq C < \infty$  where  $\|\cdot\|$  is the operator norm.

Then  $\eta$  satisfies assumption 2.

*Proof.* By intermediate value theorem, for any cube  $R_i$ ,  $1 \leq i \leq 2^{dm}$  there exists  $x_0 \in R_i$  such that  $\bar{\eta}_m(x) = \eta(x_0)$ ,  $x \in R_i$ . This implies

$$\begin{aligned} |\eta(x) - \bar{\eta}_m(x)| &= |\eta(x) - \eta(x_0)| = |\nabla \eta(\xi) \cdot (x - x_0)| \leq \\ &\leq \|\nabla \eta(\xi)\|_1 \|x - x_0\|_\infty \leq S \cdot 2^{-m} \end{aligned}$$

On the other hand, if  $R_i \subset A_{t_*}$  then

$$\begin{aligned} |\eta(x) - \bar{\eta}_m(x)| &= |\eta(x) - \eta(x_0)| = \\ &= |\nabla \eta(x_0) \cdot (x - x_0) + \frac{1}{2}[D^2 \eta(\xi)](x - x_0) \cdot (x - x_0)| \geq \\ &\geq |\nabla \eta(x_0) \cdot (x - x_0)| - \frac{1}{2} \sup_{\xi} \|D^2 \eta(\xi)\| \max_{x \in R_i} \|x - x_0\|_2^2 \geq \\ &\geq |\nabla \eta(x_0) \cdot (x - x_0)| - C_1 2^{-2m} \end{aligned} \tag{A.2}$$

Note that a strictly positive continuous function

$$h(y, u) = \int_{[0,1]^d} (u \cdot (x - y))^2 dx$$

achieves its minimal value  $h_* > 0$  on a compact set  $[0, 1]^d \times \{u \in \mathbb{R}^d : \|u\|_1 = 1\}$ .

This implies (using (A.2) and the inequality  $(a - b)^2 \geq \frac{a^2}{2} - b^2$ )

$$\begin{aligned} \Pi^{-1}(R_i) \int_{R_i} (\eta(x) - \bar{\eta}_m(x))^2 p(x) dx &\geq \\ &\geq \frac{1}{2} (p_2 2^{dm})^{-1} \int_{R_i} (\nabla \eta(x_0) \cdot (x - x_0))^2 p_1 dx - C_1^2 2^{-4m} \geq \\ &\geq \frac{1}{2} \frac{p_1}{p_2} \|\nabla \eta(x_0)\|_1^2 2^{-2m} \cdot h_* - C_1^2 2^{-4m} \geq c_2 2^{-2m} \quad \text{for } m \geq m_0. \end{aligned}$$

Now take a set  $A \in \sigma(\mathcal{F}_m)$ ,  $m \geq m_0$  from assumption 2. There are 2 possibilities: either  $A \subset A_{t_*}$  or  $A \supset A_{t_*/3}$ . In the first case the computation above implies

$$\begin{aligned} \int_{[0,1]^d} (\eta - \bar{\eta}_m)^2 \Pi(dx|x \in A) &\geq c_2 2^{-2m} = \frac{c_2}{S^2} S^2 2^{-2m} \geq \\ &\geq \frac{c_2}{S^2} \|\eta - \bar{\eta}_m\|_{\infty, A}^2 \end{aligned}$$



If the second case occurs, note that, since  $\{x : 0 < |\eta(x)| < \frac{t_*}{3}\}$  has nonempty interior, it must contain a dyadic cube  $R_*$  with edge length  $2^{-m_*}$ . Then for any  $m \geq \max(m_0, m_*)$

$$\begin{aligned} & \int_{[0,1]^d} (\eta - \bar{\eta}_m)^2 \Pi(dx | x \in A) \geq \\ & \geq \Pi^{-1}(A) \int_{R_*} (\eta - \bar{\eta}_m)^2 \Pi(dx) \geq \frac{c_2}{4} 2^{-2m} \Pi(R_*) \geq \\ & \geq \frac{c_2}{S^2} \Pi(R_*) \|\eta - \bar{\eta}_m\|_{\infty, A}^2 \end{aligned}$$

and the claim follows.  $\square$

The next proposition describes conditions which allow functions to have vanishing gradient on decision boundary but requires convexity and regular behaviour of the gradient.

Everywhere below,  $\nabla\eta$  denotes the subgradient of a convex function  $\eta$ .

For  $0 < t_1 < t_2$ , define  $G(t_1, t_2) := \frac{\sup_{x \in A_{t_2} \setminus A_{t_1}} \|\nabla\eta(x)\|_1}{\inf_{x \in A_{t_2} \setminus A_{t_1}} \|\nabla\eta(x)\|_1}$ . In case when  $\nabla\eta(x)$

is not unique, we choose a representative that makes  $G(t_1, t_2)$  as small as possible.

**Proposition A.2.** *Suppose  $\eta(x)$  is Lipschitz continuous with Lipschitz constant  $S$ . Moreover, assume that there exists  $t_* > 0$  and  $q : (0, \infty) \mapsto (0, \infty)$  such that  $A_{t_*} \subset (0, 1)^d$  and*

- (a)  $b_1 t^\gamma \leq \Pi(A_t) \leq b_2 t^\gamma \ \forall t < t_*$ ;
- (b) For all  $0 < t_1 < t_2 \leq t_*$ ,  $G(t_1, t_2) \leq q\left(\frac{t_2}{t_1}\right)$ ;
- (c) Restriction of  $\eta$  to any convex subset of  $A_{t_*}$  is convex.

Then  $\eta$  satisfies assumption 2.

**Remark** The statement remains valid if we replace  $\eta$  by  $|\eta|$  in (c).

*Proof.* Assume that for some  $t \leq t_*$  and  $k > 0$

$$R \subset A_t \setminus A_{t/k}$$

is a dyadic cube with edge length  $2^{-m}$  and let  $x_0$  be such that  $\bar{\eta}_m(x) = \eta(x_0)$ ,  $x \in R$ . Note that  $\eta$  is convex on  $R$  due to (c). Using the subgradient

inequality  $\eta(x) - \eta(x_0) \geq \nabla\eta(x_0) \cdot (x - x_0)$ , we obtain

$$\begin{aligned} \int_R (\eta(x) - \eta(x_0))^2 d\Pi(x) &\geq \int_R (\eta(x) - \eta(x_0))^2 \mathcal{I} \{ \nabla\eta(x_0) \cdot (x - x_0) \geq 0 \} d\Pi(x) \\ &\geq \int_R (\nabla\eta(x_0) \cdot (x - x_0))^2 \mathcal{I} \{ \nabla\eta(x_0) \cdot (x - x_0) \geq 0 \} d\Pi(x) \end{aligned} \quad (\text{A.3})$$

The next step is to show that under our assumptions  $x_0$  can be chosen such that

$$\text{dist}_\infty(x_0, \partial R) \geq \nu 2^{-m} \quad (\text{A.4})$$

where  $\nu = \nu(k)$  is independent of  $m$ . In this case any part of  $R$  cut by a hyperplane through  $x_0$  contains half of a ball  $B(x_0, r_0)$  of radius  $r_0 = \nu(k)2^{-m}$  and the last integral in (A.3) can be further bounded below to get

$$\begin{aligned} \int_R (\eta(x) - \eta(x_0))^2 d\Pi(x) &\geq \frac{1}{2} \int_{B(x_0, r_0)} (\nabla\eta(x_0) \cdot (x - x_0))^2 p_1 dx \geq \\ &\geq c(k) \|\nabla\eta(x_0)\|_1^2 2^{-2m} 2^{-dm} \end{aligned} \quad (\text{A.5})$$

It remains to show (A.4). Assume that for all  $y$  such that  $\eta(y) = \eta(x_0)$  we have

$$\text{dist}_\infty(y, \partial R) \leq \delta 2^{-m}$$

for some  $\delta > 0$ . This implies that the boundary of the convex set

$$\{x \in R : \eta(x) \leq \eta(x_0)\}$$

is contained in  $R_\delta := \{x \in R : \text{dist}_\infty(x, \partial R) \leq \delta 2^{-m}\}$ . There are two possibilities: either  $\{x \in R : \eta(x) \leq \eta(x_0)\} \supseteq R \setminus R_\delta$  or  $\{x \in R : \eta(x) \leq \eta(x_0)\} \subset R_\delta$ .

We consider the first case only (the proof in the second case is similar). First, note that by (b) for all  $x \in R_\delta$   $\|\nabla\eta(x)\|_1 \leq q(k)\|\nabla\eta(x_0)\|_1$  and

$$\begin{aligned} \eta(x) &\leq \eta(x_0) + \|\nabla\eta(x)\|_1 \delta 2^{-m} \leq \\ &\leq \eta(x_0) + q(k)\|\nabla\eta(x_0)\|_1 \delta 2^{-m} \end{aligned} \quad (\text{A.6})$$

Let  $x_c$  be the center of the cube  $R$  and  $u$  - the unit vector in direction  $\nabla\eta(x_c)$ . Observe that

$$\begin{aligned} \eta(x_c + (1 - 3\delta)2^{-m}u) - \eta(x_c) &\geq \nabla\eta(x_c) \cdot (1 - 3\delta)2^{-m}u = \\ &= (1 - 3\delta)2^{-m} \|\nabla\eta(x_c)\|_2 \end{aligned}$$

On the other hand,  $x_c + (1 - 3\delta)2^{-m}u \in R \setminus R_\delta$  and

$$\eta(x_c + (1 - 3\delta)2^{-m}u) \leq \eta(x_0),$$

hence  $\eta(x_c) \leq \eta(x_0) - c(1 - 3\delta)2^{-m}\|\nabla\eta(x_c)\|_1$ . Consequently, for all

$$x \in B(x_c, \delta) := \left\{ x : \|x - x_c\|_\infty \leq \frac{1}{2}c2^{-m}(1 - 3\delta) \right\}$$

we have

$$\begin{aligned} \eta(x) &\leq \eta(x_c) + \|\nabla\eta(x_c)\|_1\|x - x_c\|_\infty \leq \\ &\leq \eta(x_0) - \frac{1}{2}c2^{-m}(1 - 3\delta)\|\nabla\eta(x_c)\|_1 \end{aligned} \quad (\text{A.7})$$

Finally, recall that  $\eta(x_0)$  is the average value of  $\eta$  on  $R$ . Together with (A.6), (A.7) this gives

$$\begin{aligned} \Pi(R)\eta(x_0) &= \int_R \eta(x)d\Pi = \int_{R_\delta} \eta(x)d\Pi + \int_{R \setminus R_\delta} \eta(x)d\Pi \leq \\ &\leq (\eta(x_0) + q(k)\|\nabla\eta(x_0)\|_1\delta 2^{-m})\Pi(R_\delta) + \\ &+ (\eta(x_0) - c_2 2^{-m}(1 - 3\delta)\|\nabla\eta(x_0)\|_1)\Pi(B(x_c, \delta)) + \\ &+ \eta(x_0)\Pi(R \setminus (R_\delta \cup B(x_c, \delta))) = \\ &= \Pi(R)\eta(x_0) + q(k)\|\nabla\eta(x_0)\|_1\delta 2^{-m}\Pi(R_\delta) - \\ &- c_2 2^{-m}(1 - 3\delta)\|\nabla\eta(x_0)\|_1\Pi(B(x_c, \delta)) \end{aligned}$$

Since  $\Pi(R_\delta) \leq p_2 2^{-dm}$  and  $\Pi(B(x_c, \delta)) \geq c_3 2^{-dm}(1 - 3\delta)^d$ , the inequality above implies

$$c_4 q(k)\delta \geq (1 - 3\delta)^{d+1}$$

which is impossible for small  $\delta$  (e.g., for  $\delta < \frac{c}{q(k)(3d+4)}$ ).

Let  $A$  be a set from condition 2. If  $A \supseteq A_{t_*/3}$ , then there exists a dyadic cube  $R_*$  with edge length  $2^{-m_*}$  such that  $R_* \subset A_{t_*/3} \setminus A_{t_*/k}$  for some  $k > 0$ , and the claim follows from (A.5) as in proposition A.1.

Assume now that  $A_t \subset A \subset A_{3t}$  and  $3t \leq t_*$ . Condition (a) of the proposition implies that for any  $\varepsilon > 0$  we can choose  $k(\varepsilon) > 0$  large enough so that

$$\Pi(A \setminus A_{t/k}) \geq \Pi(A) - b_2(t/k)^\gamma \geq \Pi(A) - \frac{b_2}{b_1}k^{-\gamma}\Pi(A_t) \geq (1 - \varepsilon)\Pi(A) \quad (\text{A.8})$$

This means that for any partition of  $A$  into dyadic cubes  $R_i$  with edge length  $2^{-m}$  at least half of them satisfy

$$\Pi(R_i \setminus A_{t/k}) \geq (1 - c\varepsilon)\Pi(R_i) \quad (\text{A.9})$$

Let  $\mathcal{I}$  be the index set of cardinality  $|\mathcal{I}| \geq c\Pi(A)2^{dm-1}$  such that (A.9) is true for  $i \in \mathcal{I}$ . Since  $R_i \cap A_{t/k}$  is convex, there exists<sup>3</sup>  $z = z(\varepsilon) \in \mathbb{N}$  such that for any such cube  $R_i$  there exists a dyadic sub-cube with edge length  $2^{-(m+z)}$  entirely contained in  $R_i \setminus A_{t/k}$ :

$$T_i \subset R_i \setminus A_{t/k} \subset A_{3t} \setminus A_{t/k}.$$

It follows that  $\Pi(\bigcup_i T_i) \geq \tilde{c}(\varepsilon)\Pi(A)$ . Recall that condition (b) implies

$$\frac{\sup_{x \in \bigcup_i T_i} \|\nabla \eta(x)\|_1}{\inf_{x \in \bigcup_i T_i} \|\nabla \eta(x)\|_1} \leq q(3k)$$

Finally,  $\sup_{x \in A_{3t}} \|\nabla \eta(x)\|_2$  is attained at the boundary point, that is for some  $x_* : |\eta(x_*)| = 3t$ , and by (b)

$$\sup_{x \in A_{3t}} \|\nabla \eta(x)\|_1 \leq \sqrt{d} \|\nabla \eta(x_*)\|_1 \leq q(3k) \sqrt{d} \inf_{x \in A_{3t} \setminus A_{t/k}} \|\nabla \eta(x)\|_1.$$

Application of (A.5) to every cube  $T_i$  gives

$$\begin{aligned} \sum_{i \in \mathcal{I}} \int_{T_i} (\eta(x) - \bar{\eta}_{m+z}(x))^2 d\Pi(x) &\geq c_1(k) \Pi(A) |\mathcal{I}| \inf_{x \in A_{3t} \setminus A_{t/k}} \|\nabla \eta(x)\|_1^2 2^{-2m} 2^{-dm} \geq \\ &\geq c_2(k) \Pi(A) \sup_{x \in A_{3t}} \|\nabla \eta(x)\|_1^2 2^{-2m} \geq c_3(k) \Pi(A) \|\eta - \bar{\eta}(m)\|_{\infty, A}^2 \end{aligned}$$

concluding the proof.  $\square$

## Appendix B: Proof of Theorem 3.3

The main ideas of this proof, which significantly simplifies and clarifies initial author's version, are due to V. Koltchinskii. For convenience and brevity, let us introduce additional notations. Recall that

$$s_m = m(s + \log \log_2 N)$$

---

<sup>3</sup>If, on the contrary, every sub-cube with edge length  $2^{-(m+z)}$  contains a point from  $A_{t/k}$ , then  $A_{t/k}$  must contain the convex hull of these points which would contradict (A.8) for large  $z$ .

Let

$$\begin{aligned}\tau_N(m, s) &:= K_1 \frac{2^{dm} + s_m}{N} \\ \pi_N(m, s) &:= K_2 \frac{2^{dm} + s + \log \log_2 N}{N}\end{aligned}$$

By  $\mathcal{E}_P(\mathcal{F}, f)$  (or  $\mathcal{E}_{P_N}(\mathcal{F}, f)$ ) we denote the excess risk of  $f \in \mathcal{F}$  with respect to the true (or empirical) measure:

$$\begin{aligned}\mathcal{E}_P(\mathcal{F}, f) &:= P(y - f(x))^2 - \inf_{g \in \mathcal{F}} P(y - g(x))^2 \\ \mathcal{E}_{P_N}(\mathcal{F}, f) &:= P_N(y - f(x))^2 - \inf_{g \in \mathcal{F}} P_N(y - g(x))^2\end{aligned}$$

It follows from Theorem 4.2 in Koltchinskii [10] and the union bound that there exists an event  $\mathcal{B}$  of probability  $\geq 1 - e^{-s}$  such that on this event the following holds for all  $m$  such that  $dm \leq \log N$ :

$$\begin{aligned}\mathcal{E}_P(\mathcal{F}_m, \hat{f}_{\hat{m}}) &\leq \pi_N(m, s) \\ \forall f \in \mathcal{F}_m, \quad \mathcal{E}_P(\mathcal{F}_m, f) &\leq 2(\mathcal{E}_{P_N}(\mathcal{F}_m, f) \vee \pi_N(m, s)) \\ \forall f \in \mathcal{F}_m, \quad \mathcal{E}_{P_N}(\mathcal{F}_m, f) &\leq \frac{3}{2}(\mathcal{E}_P(\mathcal{F}_m, f) \vee \pi_N(m, s)).\end{aligned}\tag{B.1}$$

We will show that on  $\mathcal{B}$ ,  $\{\hat{m} \leq \bar{m}\}$  holds. Indeed, assume that, on the contrary,  $\hat{m} > \bar{m}$ ; by definition of  $\hat{m}$ , we have

$$P_N(Y - \hat{f}_{\hat{m}})^2 + \tau_N(\hat{m}, s) \leq P_N(Y - \hat{f}_{\bar{m}})^2 + \tau_N(\bar{m}, s),$$

which implies

$$\mathcal{E}_{P_N}(\mathcal{F}_{\hat{m}}, \hat{f}_{\bar{m}}) \geq \tau_N(\hat{m}, s) - \tau_N(\bar{m}, s) > 3\pi_N(\hat{m}, s)$$

for  $K_1$  big enough. By (B.1),

$$\mathcal{E}_{P_N}(\mathcal{F}_{\hat{m}}, \hat{f}_{\bar{m}}) = \inf_{f \in \mathcal{F}_{\bar{m}}} \mathcal{E}_{P_N}(\mathcal{F}_{\hat{m}}, f) \leq \frac{3}{2} \left( \inf_{f \in \mathcal{F}_{\bar{m}}} \mathcal{E}_P(\mathcal{F}_{\hat{m}}, f) \vee \pi_N(\hat{m}, s) \right),$$

and combination the two inequalities above yields

$$\inf_{f \in \mathcal{F}_{\bar{m}}} \mathcal{E}_P(\mathcal{F}_{\hat{m}}, f) > \pi_N(\hat{m}, s)\tag{B.2}$$

Since for any  $m$   $\mathcal{E}_P(\mathcal{F}_m, f) \leq \mathbb{E}(f(X) - \eta(X))^2$ , the definition of  $\bar{m}$  and (B.2) imply that

$$\pi_N(\bar{m}, s) \geq \inf_{f \in \mathcal{F}_{\bar{m}}} \mathbb{E}(f(X) - \eta(X))^2 > \pi_N(\hat{m}, s),$$

contradicting our assumption, hence proving the claim.

## References

- [1] J.-Y. Audibert and A. B. Tsybakov. Fast learning rates for plug-in classifiers. *Preprint*, 2005. Available at: [http://imagine.enpc.fr/publications/papers/05preprint\\_AudTsy.pdf](http://imagine.enpc.fr/publications/papers/05preprint_AudTsy.pdf).
- [2] M.-F. Balcan, S. Hanneke, and J. Wortman. The true sample complexity of active learning. In *COLT*, pages 45–56, 2008.
- [3] M.-F. Balcan, A. Beygelzimer, and J. Langford. Agnostic active learning. *J. Comput. System Sci.*, 75(1):78–89, 2009.
- [4] R. M. Castro and R. D. Nowak. Minimax bounds for active learning. *IEEE Trans. Inform. Theory*, 54(5):2339–2353, 2008.
- [5] S. Gaïffas. Sharp estimation in sup norm with random design. *Statist. Probab. Lett.*, 77(8):782–794, 2007.
- [6] E. Giné and R. Nickl. Confidence bands in density estimation. *Ann. Statist.*, 38(2):1122–1170, 2010.
- [7] S. Hanneke. Rates of convergence in active learning. *Ann. Statist.*, 39(1):333–361, 2011.
- [8] M. Hoffmann and R. Nickl. On adaptive inference and confidence bands. *The Annals of Statistics*, (to appear), 2011.
- [9] V. Koltchinskii. Rademacher complexities and bounding the excess risk in active learning. *J. Mach. Learn. Res.*, 11:2457–2485, 2010.
- [10] V. Koltchinskii. *Oracle inequalities in empirical risk minimization and sparse recovery problems*. Springer, 2011. Lectures from the 38th Probability Summer School held in Saint-Flour, 2008, École d’Été de Probabilités de Saint-Flour.
- [11] M. G. Low. On nonparametric confidence intervals. *Ann. Statist.*, 25(6):2547–2554, 1997.
- [12] A. B. Tsybakov. Optimal aggregation of classifiers in statistical learning. *Ann. Statist.*, 32(1):135–166, 2004.
- [13] A. B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer, 2009.
- [14] A. W. van der Vaart and J. A. Wellner. *Weak convergence and empirical processes*. Springer Series in Statistics.